

On Convergence Property of Implicit Self-paced Objective

Zilu Ma¹, Shiqi Liu¹ & Deyu Meng^{1*}

¹*Institute for Information and System Sciences and Ministry of Education Key Lab of Intelligent Networks and Network Security, Xi'an Jiaotong University, Xi'an 710049, China*

Abstract Self-paced learning (SPL) is a new methodology that simulates the learning principle of humans/animals to start learning easier aspects of a learning task, and then gradually take more complex examples into training. This new-coming learning regime has been empirically substantiated to be effective in various computer vision and pattern recognition tasks. Recently, it has been proved that the SPL regime has a close relationship to a implicit self-paced objective function. While this implicit objective could provide helpful interpretations to the effectiveness, especially the robustness, insights under the SPL paradigms, there are still no theoretical results strictly proved to verify such relationship. To this issue, in this paper, we provide some convergence results on this implicit objective of SPL. Specifically, we prove that the learning process of SPL always converges to critical points of this implicit objective under some mild conditions. This result verifies the intrinsic relationship between SPL and this implicit objective, and makes the previous robustness analysis on SPL complete and theoretically rational.

Keywords Self-paced learning, machine learning, non-convex optimization, convergence

Citation Zilu Ma, Shiqi Liu, Deyu Meng, et al. On Convergence Property of Implicit Self-paced Objective. , for review

1 Introduction

Self-paced learning (SPL) is a recently raised methodology designed through simulating the learning principle of humans/animals [3]. A variety of SPL realization schemes have been designed, and empirically substantiated to be effective in different computer vision and pattern recognition tasks, such as object detector adaptation [10], specific-class segmentation learning [2], visual category discovery [4], concept learning [5], long-term tracking [9], graph trimming [14], co-saliency detection [16], matrix factorization [17], face identification [6], and multimedia event detection [13].

To explain the underlying effectiveness mechanism inside SPL, [7] firstly provided some new theoretical understandings under the SPL scheme. Specifically, this work proved that the alternative optimization strategy (AOS) on SPL accords with a majorization minimization (MM) algorithm implemented on an implicit objective function. Furthermore, it is found that the loss function contained in this implicit objective has a similar configuration with non-convex regularized penalty (NCRP), leading to a rational interpretation to the robustness insight under SPL.

*Corresponding author (email: dymeng@mail.xjtu.edu.cn)

However, such understanding is still not theoretically strict. The theory in [7] can only guarantee that during the iterations of SPL solving process (i.e., the MM algorithm), the implicit objective is monotonically decreasing, while cannot prove any convergence results on this implicit objective theoretically. However, this theoretical result regarding this implicit objective is critical to the soundness of the robustness insight explanation of SPL, which guarantees to settle the convergence point of the algorithm down on the expected implicit objective, and intrinsically relate the original SPL model and this implicit objective.

To this theoretical issue of SPL, in this paper, we prove that the optimization of the implicit objective actually converges to critical points of original SPL problem under satisfactorily weak conditions. This result provides an affirmative answer to our guess that the SPL intrinsically optimizes a robust implicit objective.

In what follows, we will first introduce some related background of this research, and then we provide the main theoretical result of this work.

2 Related work

In this section, we first briefly introduce the definition of SPL, and then provides its relationship to the implicit objective of NCRP.

2.1 The SPL objective

Given training data set $\{(x_i, y_i)\}_{i=1}^N$, many machine learning problems need to minimizing the following form of objective function:

$$J(w) = \phi_\lambda(w) + \sum_{i=1}^N L(y_i, g(x_i, w)),$$

where $w \in \mathbb{R}^D$ is variables to be solved, ϕ_λ is a regularizer parameter, L is the loss function and $g(\cdot, w)$ is the parametrized learning machine, like a discriminative or a regression function.

To improve the robustness, specially avoiding the negative influence brought by large-noise-outliers, SPL imposes additional importance weights $v = (v_1, \dots, v_n)$ to loss functions of all samples, adjusted by a self-paced regularizer (SP-regularizer). Here, each $v_i \in [0, 1]$ represents how much extent the sample (x_i, y_i) will be trained in the learning process. The **self-paced objective** can then be designed as [1]:

$$E(w, v; \lambda) = \phi_\lambda(w) + \sum_{i=1}^N v_i L(y_i, g(x_i, w)) + f_\lambda(v), \quad (1)$$

where f is the **SP-regularizer**, satisfying the following conditions:

1. $v \mapsto f_\lambda(v)$ is convex on $[0, 1]$;
2. Let

$$v_\lambda^*(l) = \arg \min_{v \in [0, 1]} \{vl + f_\lambda(v)\},$$

then $l \mapsto v_\lambda^*(l)$ is non-increasing, and

$$\lim_{l \rightarrow 0} v_\lambda^*(l) = 1, \quad \lim_{l \rightarrow \infty} v_\lambda^*(l) = 0;$$

3. $\lambda \mapsto v_\lambda^*(l)$ is non-decreasing, and

$$\lim_{\lambda \rightarrow 0} v_\lambda^*(l) = 0, \quad \lim_{\lambda \rightarrow \infty} v_\lambda^*(l) \leq 1.$$

Throughout this paper, we shall assume that $v_\lambda^*(l)$ can be uniquely determined and thus can be seen as a real-valued function instead of a set-valued function.

The three conditions in the definition above provide basic principles for constructing a SP-regularizer. Condition 2 indicates that the model inclines to select easy samples (with smaller losses) in favor of complex samples (with larger losses). Condition 3 states that when the model “pace” (controlled by the pace parameter λ) gets larger, it tends to incorporate more, probably complex, samples to train a “mature” model. The convexity in Condition 1 further ensures the soundness of this regularizer for optimization.

The existence of the SP-regularizer can be illustrated by the following example.

Let the SP-regularizer be

$$f_\lambda(v) = \lambda v(\log v - 1),$$

then it yields

$$v_\lambda^*(l) = e^{-\lambda^{-1}l}.$$

It is easy to verify that $v_\lambda^*(l)$ satisfies the above conditions.

In the following, we shall write:

$$l_i(w) = L(y_i, g(x_i, w)), \quad i = 1, \dots, N$$

for simplicity.

2.2 The implicit NCRP objective

Let

$$F_\lambda(l) = \int_0^l v_\lambda^*(\tau) d\tau.$$

Since v_λ^* is non-increasing, the set of its discontinuous points is countable and consists only of jump discontinuity. Thus v_λ^* is integrable and F_λ is absolutely continuous and concave. We now define

$$G_\lambda(w) = \phi_\lambda(w) + \sum_{i=1}^N (F_\lambda \circ l_i)(w) \quad (2)$$

as the **implicit objective**, where $g \circ f$ denotes that g composed with f . An interesting observation is that this implicit SPL objective has a close relationship to NCRP widely investigated in machine learning and statistics, which provides some helpful explanation to the robustness insight under SPL [7].

The original utilized AOS algorithm for solving the SPL problem is designed by performing coordinate descent calculation on $E(w, v; \lambda)$, i.e., iterating through the process as:

$$(w^{k-1}, v^{k-1}) \rightarrow (w^{k-1}, v^k) \rightarrow (w^k, v^k).$$

Specifically, given (w^0, v^0) , if we have finished $(k-1)$ steps, then the AOS algorithm need to iteratively calculating the following two subproblems:

$$\begin{aligned} v^k &= \arg \min_v E(w^{k-1}, v; \lambda) = \arg \min_v \left\{ \sum_{i=1}^N v_i l_i(w^{k-1}) + f_\lambda(v_i) \right\}, \\ w^k &\in \arg \min_w E(w, v^k; \lambda) = \arg \min_w \left\{ \phi_\lambda(w) + \sum_{i=1}^N v_i^k l_i(w) \right\}. \end{aligned}$$

Note that the first subproblem is feasible since we have assumed that v_λ^* can be uniquely determined. Indeed, using the notation of v_λ^* , we have

$$v_i^k = v_\lambda^*(l_i(w^{k-1})), \quad i = 1, \dots, N.$$

We then set

$$Q(w|w^*) = \sum_{i=1}^N (F_\lambda \circ l_i)(w^*) + (v_\lambda^* \circ l_i)(w^*)[l_i(w) - l_i(w^*)].$$

It is easy to deduce that $Q(w|w^*)$ is actually the first-order Taylor series of F_λ at $l_i(w^*)$. Based on the concavity of F_λ , we know that

$$U(w|w^*) = \phi_\lambda(w) + Q(w|w^*)$$

constitutes a upper bound of $G_\lambda(w)$ (as defined in Eq. (2)), which provides a qualified surrogate function for MM algorithm.

One of the key issues in [7] is that if $\{w^k\}$ is produced by AOS algorithm of $E(w, v; \lambda)$, then it can also be produced by performing MM algorithm on G_λ and vice versa. We prove one side by induction. The other side is totally the same. Suppose we have proved that w^k can be produced by performing MM algorithm on G_λ at k^{th} step. When it comes to the $(k+1)^{\text{th}}$ step,

$$\begin{aligned} w^{k+1} &\in \arg \min_w E(w, v^{k+1}; \lambda) \\ &= \arg \min_w \phi_\lambda(w) + \sum_{i=1}^N v_i^{k+1} l_i(w) \\ &= \arg \min_w \phi_\lambda(w) + \sum_{i=1}^N v_\lambda^*(l_i(w^k)) \cdot l_i(w) \\ &= \arg \min_w \phi_\lambda(w) + Q(w|w^k) = \arg \min_w U(w|w^k). \end{aligned} \quad (3)$$

Thus we have proved our claim that these two optimization algorithms (AOS/MM) conducting on the two different objective functions ($E(w, v; \lambda)/G_\lambda(w)$) are intrinsically equivalent.

We then need to prove whether every convergence point of MM algorithm, or equivalently, that of the AOS algorithm on the SPL objective, is at least a critical point of G_λ .

3 The main convergence result

Actually, the proof of the convergence of MM algorithm is basically the same as that of the EM algorithm (see [12]) only with some obvious changes, as discussed in [11]. And the convergence of EM and MM is indeed a corollary of a global convergence theorem of Zangwill (see [15]). We can generalize the proof to the case of variational analysis. Before that, we need to clarify some terminologies which can be referred to in [8].

A function $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ is said to be **lower semi-continuous** or simply **lsc** if

$$\text{lev}_{f \leq \alpha} := \{x : f(x) \leq \alpha\}$$

is closed for any $\alpha \in \mathbb{R}$. f is said to be **level-bounded** if $\text{lev}_{f \leq \alpha}$ is bounded for any α . And f is called **coercive** if $\lim_{|x| \rightarrow \infty} f(x) = \infty$. Note that coercive functions are level-bounded. A **critical point** x of f means that $0 \in \partial f(x)$, where ∂ stands for the **subdifferential** [8].

The main theorem of this paper can then be stated as follows.

Theorem 1. Suppose that the objective function of MM algorithm, $G : \mathbb{R}^D \rightarrow \overline{\mathbb{R}}$, is lsc and level-bounded, and that the surrogate function at w^* is $U(\cdot|w^*)$, which is lsc as a function on \mathbb{R}^{2D} , and satisfies

$$\partial U(w|w) \subset \partial G(w), \quad \forall w \in \mathbb{R}^D,$$

where $\partial U(w|w^*)$ is the partial subdifferential in w . Then for any initial parameter w^0 , every cluster point of the produced sequence $\{w^k\}$ of MM algorithm is a critical point of G .

Proof. See the appendix.

For our problem, we can give a sufficient condition of convergence, which is easy to verify and satisfied by most of the current SPL variations.

Theorem 2. In the SPL objective as defined in Section 2.1, suppose L is bounded below, $w \mapsto L(y, g(x, w))$ is continuously differentiable, $v_\lambda^*(\cdot)$ is continuous, and ϕ_λ is coercive and lsc. Then for any initial parameter w^0 , every cluster point of the produced sequence $\{w^k\}$, obtained by the AOS algorithm on solving Eq. (1), is a critical point of the implicit objective G_λ as defined in Eq. (2).

Proof. It is obvious that G_λ is lsc and level-bounded and U is lsc as a function on \mathbb{R}^{2D} with these assumptions. And the continuity of v_λ^* makes F_λ continuously differentiable. Then we have

$$\begin{aligned} \partial G_\lambda(w^*) &= \partial \phi_\lambda(w^*) + \sum_{i=1}^N F'_\lambda(l_i(w^*)) \nabla l_i(w^*) \\ &= \partial \phi_\lambda(w^*) + \sum_{i=1}^N (v_\lambda^* \circ l_i)(w^*) \nabla l_i(w^*) \\ &= \partial U(w^* | w^*). \end{aligned}$$

Based on Theorem 1, for any initial parameter w^0 , every cluster point of the produced sequence $\{w^k\}$ is a critical point of G_λ . The proof is then completed.

From the theorem, we can see that the AOS algorithm generally used to solving the SPL problem can be guaranteed to convergent to a critical point of the implicit NCRP objective G_λ . The intrinsic relationship between two objectives can then be constructed.

Note that in the above theorem, it is required that every minimization step in MM algorithm exactly attains the minima of the surrogate function $U(w|w^k)$, i.e.,

$$U(w^{k+1} | w^k) = \min U(\cdot | w^k). \quad (4)$$

This is generally hard to achieve in real applications, especially for those learning models without closed-form solution. We thus want to further relax the condition to allow a relatively weaker solution “with errors” in implementing the MM algorithm on the surrogate function. That is, we can weaken the condition (4) as:

$$U(w^{k+1} | w^k) \leq \min U(\cdot | w^k) + \epsilon_k,$$

where $\epsilon_1, \epsilon_2, \dots$ is a non-negative sequence satisfying $\{\epsilon_k\} \in l^1$, i.e., $\sum_k \epsilon_k < \infty$.

Under this relaxed condition, we can still prove the convergence result of SPL in the following algorithm

Theorem 3. In the SPL objective as defined in Section 2.1, suppose L is bounded below, $w \mapsto L(y, g(x, w))$ is continuously differentiable, $v_\lambda^*(\cdot)$ is continuous, and ϕ_λ is coercive and lsc. Let w^0 be an arbitrary initial parameter, and $\{w^k\}$ be the sequence obtained by the AOS algorithm on solving Eq. (1) with errors $\{\epsilon_k \geq 0\} \in l^1$, that is,

$$E(w^k, v^k; \lambda) \leq \min E(\cdot, v^k; \lambda) + \epsilon_k, \quad \forall k \geq 1.$$

Then every cluster point of $\{w^k\}$ is a critical point of the implicit objective G_λ as defined in Eq. (2).

Based on the theorem, we can then confirm the intrinsic relationship between SPL and its implicit objective.

4 Conclusion

In this paper, we have proved that the learning process of traditional SPL regime can be guaranteed to converge to rational critical points of the corresponding implicit NCRP objective. This theory helps confirm the intrinsic relationship between SPL and this implicit objective, and thus verifies previous robustness analysis of SPL on the basis of the understanding of such relationship. Besides, we have used

some new theoretical skills for the proof of convergence, which inclines to be beneficial to the previous MM and EM convergence theories to a certain extent.

Conflict of interest The authors declare that they have no conflict of interest.

References

- 1 L. Jiang, D. Meng, T. Mitamura, and A. Hauptmann. Easy samples first: self-paced reranking for zeroexample multimedia search. In *ACM MM*, 2014.
- 2 M. Kumar, H. Turki, D. Preston, and D. Koller. Learning specific-class segmentation from diverse data. In *ICCV*, 2011.
- 3 M Pawan Kumar, Benjamin Packer, and Daphne Koller. Self-paced learning for latent variable models. In *Advances in Neural Information Processing Systems*, pages 1189–1197, 2010.
- 4 Y. Lee and K. Grauman. Learning the easy things first: Self-paced visual category discovery. In *CVPR*, 2011.
- 5 Junwei Liang, Lu Jiang, Deyu Meng, and Alex Hauptmann. Learning to detect concepts from webly-labeled video data. In *IJCAI*, 2016.
- 6 L. Lin, K. Wang, D. Meng, W. Zuo, and L. Zhang. Active self-paced learning for cost-effective and progressive face identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- 7 D. Meng, Q. Zhao, and L. Jiang. What Objective Does Self-paced Learning Indeed Optimize? *ArXiv e-prints*, November 2015.
- 8 R. Tyrrell Rockafellar and Roger J-B Wets. *Variational analysis*, volume 317. Springer Science & Business Media, 2009.
- 9 J. Supančič III and D. Ramanan. Self-paced learning for long-term tracking. In *CVPR*, 2013.
- 10 K. Tang, V. Ramanathan, F. Li, and D. Koller. Shifting weights: Adapting object detectors from image to video. In *NIPS*, 2012.
- 11 Florin Vaida. Parameter convergence for EM and MM algorithms. *Statistica Sinica*, pages 831–840, 2005.
- 12 CF Jeff Wu. On the convergence properties of the EM algorithm. *The Annals of statistics*, pages 95–103, 1983.
- 13 S. Yu, L. Jiang, Z. Mao, and et al. CMU-Informedia@TRECVID 2014 multimedia eventdetection (MED). In *TRECVID Video Retrieval Evaluation Workshop*, 2014.
- 14 Zongsheng Yue, Deyu Meng, Juan He, and Gemeng Zhang. Semi-supervised learning through adaptive laplacian graph trimming. *Image and Vision Computing*, 2016.
- 15 W.I. Zangwill. *Nonlinear programming: a unified approach*. Prentice-Hall international series in management. Prentice-Hall, 1969.
- 16 Dingwen Zhang, Deyu Meng, Chao Li, Lu Jiang, Qian Zhao, and Junwei Han. A self-paced multiple-instance learning framework for co-saliency detection. In *IEEE International Conference on Computer Vision*, pages 594–602, 2015.
- 17 Qian Zhao, Deyu Meng, Lu Jiang, Qi Xie, Zongben Xu, and Alexander G Hauptmann. Self-paced learning for matrix factorization. In *AAAI*, pages 3196–3202, 2015.

Appendix A Proof of Theorem 1

Theorem 1 is actually a corollary of a stronger version of Zangwill’s global convergence theorem [15, page 91]. We first need to give the following lemmas.

Lemma 1. If f is lsc, $x_n \rightarrow x$, and $\{f(x_n)\}$ is non-increasing, then $f(x_n) \rightarrow f(x)$.

Proof.

$$f(x) = \liminf_{n \rightarrow \infty} f(x_n) = \lim_{k \rightarrow \infty} \inf_{n \geq k} f(x_n) = \inf_{n \geq 1} f(x_n) = \lim_{n \rightarrow \infty} f(x_n).$$

Lemma 2. Suppose that X is a finite-dimensional Euclidean space, M is a set-valued mapping from X to X and that $\{x_k\}$ is produced by M , which means

$$x_{k+1} \in M(x_k), \quad \forall k \geq 0.$$

Γ is a subset of X that we are interested at, called the ”solution set” and satisfying

1. There is a compact subset K , such that $x_k \in K, \forall k$,
2. M is outer semicontinuous on $X \setminus \Gamma$, that is

$$x_k \rightarrow x \text{ in } X \setminus \Gamma \implies M(x_k) \rightarrow M(x).$$

3. There is a lsc function G defined on X , such that

- (a) $G(y) < G(x), \forall y \in M(x), x \notin \Gamma$,
- (b) $G(y) \leq G(x), \forall y \in M(x), x \in \Gamma$,

then all the cluster points of $\{x_k\}$ are in Γ , and $\exists \bar{x} \in \Gamma$, such that $G(x_k)$ is non-increasing and convergent to $G(\bar{x})$.

Note: we will repeatedly use the fact that $\{G(x_k)\}$ is non-increasing. Without loss of generality, we can assume that $n_1 < n_2 < \dots$ when we take a subsequence $\{x_{n_k}\}$ of $\{x_n\}$.

Proof. (1) Suppose x^* is a cluster point of $\{x_k\}$. The existence of x^* is guaranteed by the compactness of K . Thus there exists a subsequence $\{x_{n_k}\}$, such that $x_{n_k} \rightarrow x^*, (k \rightarrow \infty)$. Since $\{G(x_k)\}$ is non-increasing, based on Lemma 3, it holds that

$$G(x^*) = \lim_{k \rightarrow \infty} G(x_{n_k}).$$

Denote $G^* = G(x^*)$, and then we prove that $G(x_n) \rightarrow G^* (n \rightarrow \infty)$. This is because $\forall \epsilon > 0, \exists k_0 > 0$, such that

$$G(x_{n_k}) - G^* < \epsilon, \quad \forall k \geq k_0.$$

When $n \geq n_{k_0}$,

$$G(x_n) - G^* = G(x_n) - G(x_{n_{k_0}}) + G(x_{n_{k_0}}) - G^* < 0 + \epsilon = \epsilon.$$

There exists $k_1 > 0$, such that $n < n_{k_1}$, and thus

$$G(x_n) - G^* = G(x_n) - G(x_{n_{k_1}}) + G(x_{n_{k_1}}) - G^* \geq 0 + 0 = 0.$$

Therefore,

$$0 \leq G(x_n) - G^* < \epsilon, \quad \forall n \geq n_{k_0},$$

which indicates $G(x_n) \rightarrow G^*$.

(2) If $x^* \notin \Gamma$, take a subsequence

$$y_k = x_{n_k+1} \in M(x_{n_k}).$$

Since y_k all lie in K , there exists a subsequence $\{y_{k_l}\}$, such that $y_{k_l} \rightarrow \bar{x}, (l \rightarrow \infty)$. Since M is outer semicontinuous, $\bar{x} \in M(x^*)$. Based on Lemma 3, we know that $G(y_{k_l}) \rightarrow G(\bar{x}), (l \rightarrow \infty)$. Due to the properties of G ,

$$G(\bar{x}) < G(x^*) = \lim_{n \rightarrow \infty} G(x_n) = \lim_{l \rightarrow \infty} G(y_{k_l}) = G(\bar{x}),$$

a contradiction.

We then provide a proof of Theorem 1.

Proof. [Proof of Theorem 1]

Let Γ be the set of critical points of G , and

$$M(w^*) = \arg \min_w U(w|w^*).$$

By the descending property of MM algorithm, $G(w) \leq G(w^*), \forall w \in M(w^*)$. Condition 3b is satisfied.

Condition 1: since G is lsc and level-bounded, $K = \text{lev}_{G \leq w^0}$ is closed and bounded, and thus compact. By the descending property of MM algorithm, all the parameters w^k lie in K .

Condition 2: suppose $w^k \rightarrow w^*, v^k \rightarrow v^*, v^k \in M(w^k)$, and then $\forall w \in \mathbb{R}^D$, it holds that

$$U(v^k|w^k) \leq U(w|w^k).$$

Taking infimal limit on both sides when $k \rightarrow \infty$, we have

$$U(v^*|w^*) = \liminf_{k \rightarrow \infty} U(v^k|w^k) \leq \liminf_{k \rightarrow \infty} U(w|w^k) = U(w|w^*).$$

Thus $v^* \in M(w^*)$, which means M is outer semi-continuous.

Condition 3a: If $w^* \notin \Gamma$, then

$$0 \notin \partial G(w^*) \supset \partial U(w^*|w^*).$$

By the generalized Fermat theorem (see [8, 10.1]), w^* is not a minima of $U(\cdot|w^*)$, i.e., $w^* \notin M(w^*)$. Since $\forall w \in M(w^*)$,

$$G(w) \leq U(w|w^*) < U(w^*|w^*) = G(w^*).$$

All the conditions of the proceeding theorem are satisfied. The proof is then completed.

Appendix B Proof of Theorem 3

Similar to [8, 5.41], we give the following definition:

Definition 1. A sequence of set-valued mappings M_k **converges outer semicontinuously** to another set-valued mapping M , if

$$\limsup_{k \rightarrow \infty} M_k(x_k) \subset M(\bar{x}), \quad \forall x_k \rightarrow \bar{x},$$

that is,

$$x_k \rightarrow \bar{x}, v_k \in M_k(x_k), v_k \rightarrow \bar{v} \implies \bar{v} \in M(\bar{x}).$$

Before giving proof of Theorem 3, we need to prove the following two lemmas.

Lemma 3. Let X be an Euclidean space with finite dimension, and $M, M_k, k = 1, 2, \dots$ be set-valued mappings from X to itself. Suppose that M_k converges outer semicontinuously to M , and that $\{x_k\}$ is produced by $\{M_k\}$, which means

$$x_{k+1} \in M_k(x_k), \quad \forall k.$$

Let Γ be an arbitrary set, called the "solution set", satisfying

1. There is a compact set K such that $x_k \in K, \forall k$,
2. There is a lsc α defined on X , such that
 - (a) $\alpha(y) < \alpha(x), \forall y \in M(x), x \notin \Gamma$;
 - (b) There is a sequence of non-negative numbers $\{\epsilon_k\} \in l^1$, that is $\sum_k \epsilon_k < \infty$, and

$$\alpha(y_{k+1}) \leq \alpha(x) + \epsilon_k, \quad \forall y_{k+1} \in M_k(x), \forall x, \forall k.$$

Then all the cluster points of $\{x_k\}$ lie in Γ , and $\exists \bar{x} \in \Gamma$, such that $\alpha(x_k)$ converges to $\alpha(\bar{x})$.

Proof. (1) Set $r_k = \sum_{j \geq k} \epsilon_j$, then $r_k \rightarrow 0$, and

$$\alpha(x_{k+1}) + r_{k+1} \leq \alpha(x_k) + \epsilon_k + r_{k+1} = \alpha(x_k) + r_k.$$

Thus $\{\alpha(x_k) + r_k\}$ is non-increasing.

(2) Let x^* be a cluster point of $\{x_k\}$, and then there exists a subsequence $\{x_{n_k}\}$, such that $x_{n_k} \rightarrow x^*$. Since α is lsc, we have

$$\alpha(x^*) = \liminf_k \alpha(x_{n_k}) = \liminf_k (\alpha(x_{n_k}) + r_{n_k}) = \lim_k (\alpha(x_{n_k}) + r_{n_k}) = \lim_k \alpha(x_{n_k}).$$

The second equality holds because

$$\liminf_k \alpha(x_{n_k}) \leq \liminf_k (\alpha(x_{n_k}) + r_{n_k}) \leq \liminf_k \alpha(x_{n_k}) + \limsup_k r_{n_k} = \liminf_k \alpha(x_{n_k}).$$

And we can prove in the same way as in (1) of Lemma 2 that

$$\lim_n \alpha(x_n) = \alpha(x^*).$$

(3) We need to show that $x^* \in \Gamma$. Suppose not, take

$$y_k = x_{n_k+1} \in M_{n_k}(x_{n_k}),$$

Due to the compactness of K , there is a subsequence $\{y_{k_l}\}$ of $\{y_k\}$, such that $\exists \bar{x} \in K, y_{k_l} \rightarrow \bar{x}$. We can argue in the same way as in (2) to show that $\alpha(y_{k_l}) \rightarrow \alpha(\bar{x})$. Since

$$y_{k_l} = x_{n_{k_l}+1} \in M_{n_{k_l}}(x_{n_{k_l}}),$$

and M_k converges outer semicontinuously to M , we have $\bar{x} \in M(x^*)$. Thus

$$\alpha(\bar{x}) < \alpha(x^*) = \lim_n \alpha(x_n) = \lim_l \alpha(y_{k_l}) = \alpha(\bar{x}),$$

a contradiction.

The proof is then completed.

Now we can prove another lemma using above theoretical result.

Lemma 4. Let $F : \mathbb{R}^D \rightarrow \overline{\mathbb{R}}$ be the objective of MM algorithm. Suppose that F is lsc and level-bounded, and that the surrogate function at w^* is $U(\cdot|w^*)$. In addition, suppose $U(\cdot|w^*)$ is lsc as a function defined on \mathbb{R}^{2D} whose subgradient satisfies

$$\partial U(w|w) \subset \partial F(w), \quad \forall w \in \mathbb{R}^D,$$

where $\partial U(w|w^*)$ is the partial subdifferential with respect to w . Then for any initial parameter w^0 , all the cluster points of the sequence $\{w^k\}$ produced by MM algorithm "with errors" are still critical points of F .

Proof. We prove by a direct application of Lemma 3. Let Γ be the set consisting of all the critical points of F , $\alpha = F$, and

$$M_k(w^*) = \{w : U(w|w^*) \leq \min U(\cdot|w^*) + \epsilon_k\}.$$

M is the same as before:

$$M(w^*) = \arg \min_w U(w|w^*).$$

We first need to show that M_k converges outer semicontinuously to M . Suppose $w^k \rightarrow \bar{w}, v^k \in M_k(w^k), v^k \rightarrow \bar{v}$, and then $\forall w$,

$$U(v^k|w^k) \leq \min U(\cdot|w^k) + \epsilon_k \leq U(w|w^k) + \epsilon_k.$$

Taking infimal limit on both sides when $k \rightarrow \infty$, we have

$$U(\bar{v}|\bar{w}) = \liminf_k U(v^k|w^k) \leq \liminf_k (U(w|w^k) + \epsilon_k) \leq \liminf_k U(w|w^k) + \limsup_k \epsilon_k = U(w|\bar{w}),$$

which means $\bar{v} \in M(\bar{w})$. Thus M_k converges outer semicontinuously to M .

Condition 1: F is level-bounded, thus

$$K(w^0) = \{w : F(w) \leq F(w^0) + \sum_k \epsilon_k\}$$

is bounded. Since F is also lsc, $K(w^0)$ is closed and hence compact. By (1) of Lemma 3, w^k all lie in $K(w^0)$.

Condition 2a: If $w^* \notin \Gamma$, then

$$0 \notin \partial F(w^*) \supset \partial U(w^*|w^*).$$

By the generalized Fermat theorem, w^* is not a minima of $U(\cdot|w^*)$, and hence $w^* \notin M(w^*)$. It follows that $\forall w \in M(w^*)$,

$$F(w) \leq U(w|w^*) < U(w^*|w^*) = F(w^*).$$

Condition 2b: Let $v \in M_k(w)$, and then

$$U(v|w) \leq \min U(\cdot|w) + \epsilon_k.$$

Thus,

$$F(v) \leq U(v|w) \leq \min U(\cdot|w) + \epsilon_k \leq U(w|w) + \epsilon_k = F(w) + \epsilon_k.$$

Therefore, all the conditions of Lemma 3 are satisfied and we have finished the proof.

Just like the proof of Theorem 2, Theorem 3 can be easily proved by directly utilizing the results of the above Lemma 4. We omit the proof here.